CRAY
THE SUPERCOMPUTER COMPANY

# Future Peta-Exa Scale Architectures

## John Levesque
## February, 2011

# Outline

- Future architectural directions
- Cray XT6 and Cray XE6 compute blades
- Seastar and Gemini interconnects

# Today's Fastest System

| Systems | 2010 |
|---|---|
| System peak | 2 Pflop/s |
| Power | 6 MW |
| System memory | 0.3 PB |
| Node performance | 125 GF |
| Node memory BW | 25 GB/s |
| Node concurrency | 12 |
| Total Node Interconnect BW | 3.5 GB/s |
| System size (nodes) | 18,700 |
| Total concurrency | 225,000 |
| Storage | 15 PB |
| IO | 0.2 TB |
| MTTI | days |

3

# Today's Fastest System

| Systems | 2010 |
|---|---|
| System peak | 2 Pflop/s |
| Power | 6 MW |
| System memory | 0.3 PB |
| Node performance | 125 GF |
| Node memory BW | 25 GB/s |
| Node concurrency | 12 |
| Total Node Interconnect BW | 3.5 GB/s |
| System size (nodes) | 18,700 |
| Total concurrency | 225,000 |
| Storage | 15 PB |
| IO | 0.2 TB |
| MTTI | days |

Moving forward to Exascale…

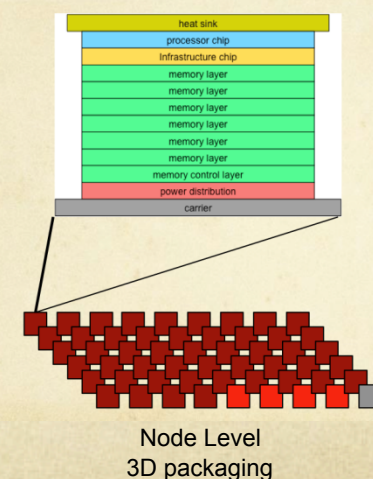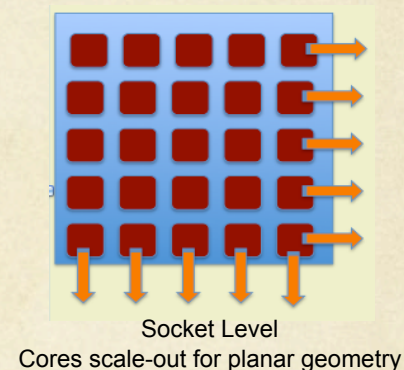The claim is that an Exascale system can cost at most $200M and consume no more than 20 MW.

4

# Potential System Architecture
## with a cap of $200M and 20MW

| Systems | 2010 | 2018 | Difference Today & 2018 |
|---|---|---|---|
| System peak | 2 Pflop/s | 1 Eflop/s | O(1000) |
| Power | 6 MW | ~20 MW | |
| System memory | 0.3 PB | 32 - 64 PB [ .03 Bytes/Flop ] | O(100) |
| Node performance | 125 GF | 1,2 or 15TF | O(10) – O(100) |
| Node memory BW | 25 GB/s [.20 Bytes/Flop] | 2 - 4TB/s [ .002 Bytes/Flop ] | O(100) |
| Node concurrency | 12 | O(1k) or 10k | O(100) – O(1000) |
| Total Node Interconnect BW | 3.5 GB/s | 200-400GB/s (1:4 or 1:8 from memory BW) | O(100) |
| System size (nodes) | 18,700 | O(100,000) or O(1M) | O(10) – O(100) |
| Total concurrency | 225,000 | O(billion) [O(10) to O(100) for latency hiding] | O(10,000) |
| Storage | 15 PB | 500-1000 PB (>10x system memory is min) | O(10) – O(100) |
| IO | 0.2 TB | 60 TB/s (how long to drain the machine) | O(100) |
| MTTI | days | O(1 day) | - O(10) |

# Exascale ($10^{18}$ Flop/s) Systems: Two possible paths

- Light weight processors (think BG/P)
  - ~1 GHz processor ($10^9$)
  - ~1 Kilo cores/socket ($10^3$)
  - ~1 Mega sockets/system ($10^6$)



Socket Level
Cores scale-out for planar geometry

- Hybrid system (think GPU based)
  - ~1 GHz processor ($10^9$)
  - ~10 Kilo FPUs/socket ($10^4$)
  - ~100 Kilo sockets/system ($10^5$)



Node Level
3D packaging

# Short Term Petascale Systems – Node Architecture

| | Cores on the node | Total threading | Vector Length | Programming Model |
|---|---|---|---|---|
| Blue Waters | 16 | 32 | 8 | OpenMP/MPI/ Vector |
| Blue Gene Q | 16 | 32 | 8 | OpenMP/MPI/ Vector |
| Magna-Cours | 24 | 24 | 4 | OpenMP/MPI/ Vector |
| OLCF3-Acc | 32 | 32 (768*) | 16 | Threads/ Cuda/Vector |
| Intel MIC | 32 | 128 | 8 | OpenMP/MPI/ Vector |
| Jaguar & Kraken | 12 | 12 | 2 | OpenMP/MPI/ Vector |

* Nvidia allows oversubscription to SIMT units

# Hybrid Multi-core Architecture




- Massively Parallel System with high powered nodes that exhibit
  - Multiple levels of parallelism
    - Shared Memory parallelism on the node
    - SIMD vector units on each core or thread
  - Potentially disparate processing units
    - Host with conventional X86 architecture
    - Accelerator with highly parallel – SIMD units
  - Potentially disparate memories
    - Host with conventional DDR memory
    - Accelerator with high bandwidth memory

- All MPI may not be best approach
  - Memory per core will decease
  - Injection bandwidth/core will decease
  - Memory bandwidth/core will decrease
- Hybrid MPI + threading on node may be able to
  - Save Memory
  - Reduce amount of off node communication required
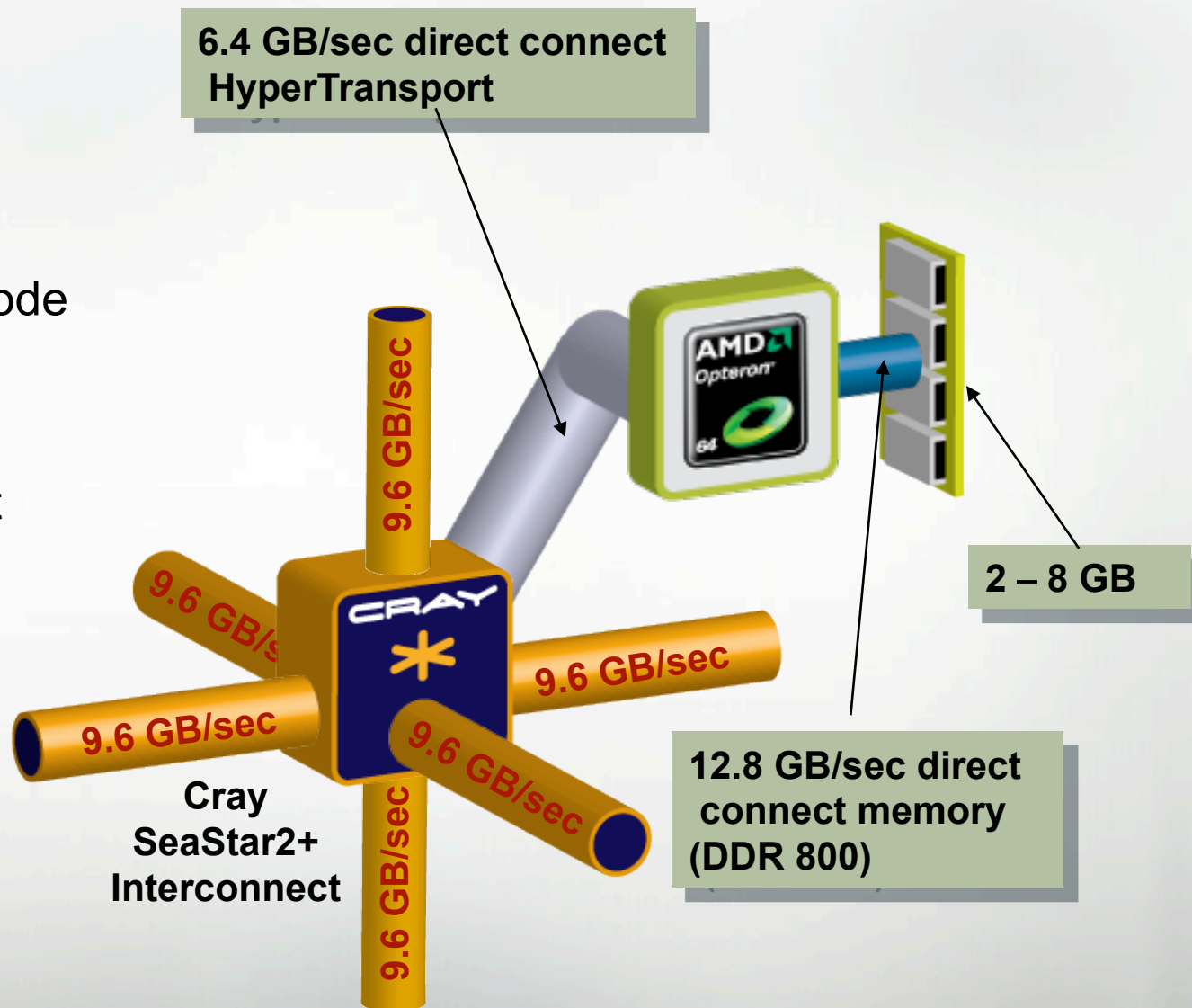  - Reduce amount of memory bandwidth required

# An overview of Cray XT systems

|  | XT3 | XT4 | XT5 | XT6 |
|---|---|---|---|---|
| Number of cores/socket | 2 | 4 | 4-6 | 12 |
| Number of cores/node | 2 | 4 | 8-12 | 24 |
| Clock Cycle (CC) | 2.6 | 2.3 | 2.6 | 2.1 |
| Number of 64 bit Results/CC | 2 | 4 | 4 | 4 |
| GFLOPS/Node | 10.4 | 36.8 | 83.6-124.8 | ~200 |
| Interconnect | Seastar 1 | Seastar 2+ | Seastar 2+ | Gemini |
| Link Bandwidth GB/sec | 6x2.4 | 6x4.8 | 6x4.8 | 10x4.7 |
| MPI Latency microseconds | 6 | 6 | 6 | 1.5 |
| Messages/sec | 400K | 400K | 400K | 10M |
| Global Addressing | No | No | No | Yes |

# Cray XT4 Node

- 4-way SMP
- >35 Gflops per node
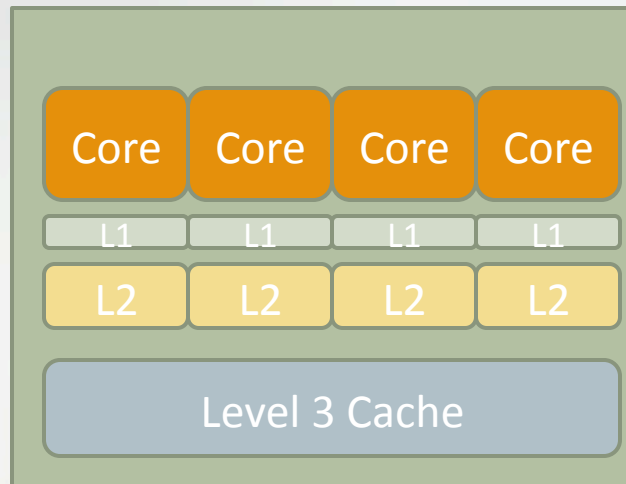- Up to 8 GB per node
- OpenMP Support within socket

6.4 GB/sec direct connect HyperTransport

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

2 – 8 GB

12.8 GB/sec direct connect memory (DDR 800)

Cray SeaStar2+ Interconnect

# Cray XT5 Node

CRAY
THE SUPERCOMPUTER COMPANY

- 8-way SMP
- >70 Gflops per node
- Up to 32 GB of shared memory per node
- OpenMP Support

2 – 32 GB memory

6.4 GB/sec direct connect HyperTransport

AMD Opteron 64

AMD Opteron 64

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

CRAY

Cray SeaStar2+ Interconnect

25.6 GB/sec direct connect memory

# XT Hardware Socket

| Core | Core | Core | Core |
|------|------|------|------|
| L1 | L1 | L1 | L1 |
| L2 | L2 | L2 | L2 |
| Level 3 Cache | | | |

| Budapest | Barcelona | Shanghai | Istanbul |
|----------|-----------|----------|----------|
| 4 FPS/CC | 4 FPS/CC | 4 FPS/CC | 4 FPS/CC |
| 64KB | 64KB | 64KB | 64KB |
| 512KB | 512KB | 512KB | 512KB |
| 2048KB | 2048KB | 6144KB | 6144KB |

# Simplified memory hierachy on the Quad Core AMD Opteron – Quad Core

registers

16 SSE2 128-bit registers
16 64 bit registers

2 x 16 Bytes per clock loads or 1 x 16 Bytes per clock store, (76.8 GB/s or 38.4 GB/s on 2.4 Ghz)

L1 data cache

16 Bytes per clock,
38.4 GB/s BW

- 64 Byte cache line
- complete data cache lines are loaded from main memory, if not in L2 or L3 cache
- if L1 data cache needs to be refilled, then storing back to L2 cache, if L2 needs to be refilled, storing back to L3

L2 cache

. . . . . .

32 GB/s

- Items in L1 and L2 are exclusive, L3 is "sharing aware"
- write back cache: data offloaded from L1 data cache are stored here first
  until they are flushed out to main memory
- Pre-fetching can go to L1 or L2 cache

Shared L3 cache

. . . . . .

16 Bytes wide => 12.8 GB/s for DDR2-800, 73ns

8GB/s over coherent Hyper Transport, 115ns

Main memory

Remote memory

# XT Hardware Node

Interconnect

1 HT Links

6.4 GB/sec

2 HT Links

12.8 GB/sec

| Core | Core | Core | Core |
|---|---|---|---|
| L1 | L1 | L1 | L1 |
| L2 | L2 | L2 | L2 |

Level 3 Cache

73 ns
12.8 GB/sec

115 ns
8 GB/sec

Memory

| Core | Core | Core | Core |
|---|---|---|---|
| L1 | L1 | L1 | L1 |
| L2 | L2 | L2 | L2 |

Level 3 Cache

73 ns

12.8 GB/sec

Memory

© Cray Inc.

- **Strengths**
  - Upgradability
  - Scalability of interconnects
  - Increased node performance – need to use fewer nodes to achieve same performance
  - Global addressing with Gemini
    - Adds ability to use PGAS – UPC And CAF to program around some of the weaknesses

## Bottlenecks

- Memory bandwidth will never be enough to support all the cores
  - ➢ Need to think about programming around this
    - o PGAS – UPC and CAF
    - o OpenMP
- Injection Bandwidth will never be enough to support all the cores
  - ➢ Need to think about programming around this
    - o PGAS – UPC and CAF
    - o OpenMP
- Global Bandwidth will never be enough to support ALL to ALLs across all of the cores
  - ➢ Need to think about programming around this
    - o PGAS – UPC and CAF
    - o OpenMP

# XE6 Node Details: 24-core Magny Cours



- 2 Multi-Chip Modules, 4 Opteron Dies
- 8 Channels of DDR3 Bandwidth to 8 DIMMs
- 24 (or 16) Computational Cores
  - 64 KB L1 and 512 KB L2 caches for each core
  - 6 MB of shared L3 cache on each die
- Dies are fully connected with HT3
- Snoop Filter Feature Allows 4 Die SMP to scale well

# XE6 Node Details: 24-core Magny Cours

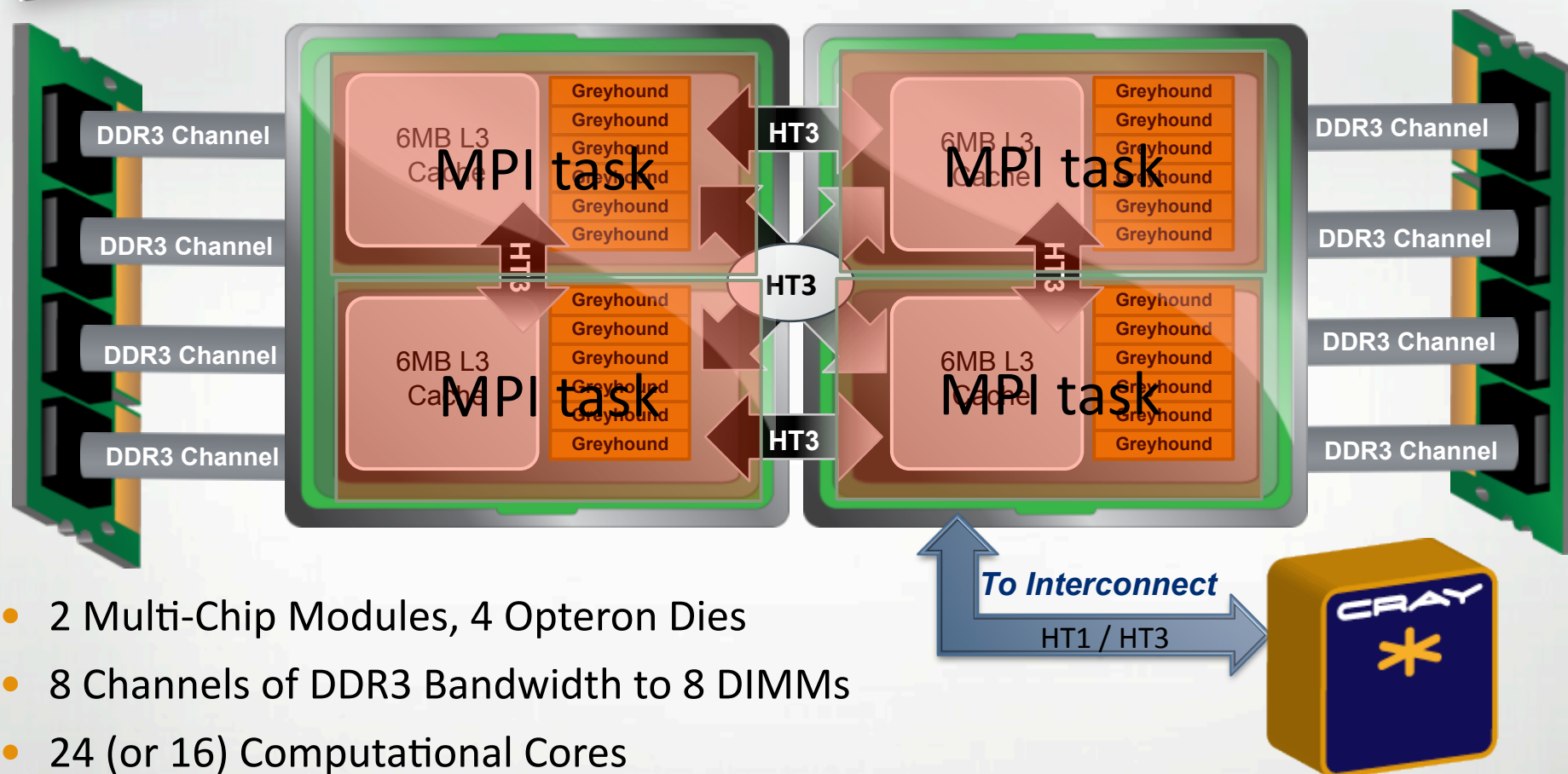Run using 2 MPI tasks on the node
One on Each Die

MPI task    MPI task

To Interconnect
HT1 / HT3

- 2 Multi-Chip Modules, 4 Opteron Dies
- 8 Channels of DDR3 Bandwidth to 8 DIMMs
- 24 (or 16) Computational Cores
  - 64 KB L1 and 512 KB L2 caches for each core
  - 6 MB of shared L3 cache on each die
- Dies are fully connected with HT3
- Snoop Filter Feature Allows 4 Die SMP to scale well

Use OpenMP across all 12 cores in the Die

# XE6 Node Details: 24-core Magny Cours

Run using 4 MPI tasks on the node
One on Each Socket

DDR3 Channel

DDR3 Channel

DDR3 Channel

DDR3 Channel

6MB L3 Cache

MPI task

Greyhound
Greyhound
Greyhound
Greyhound
Greyhound
Greyhound

HT3

HT3

6MB L3 Cache

MPI task

Greyhound
Greyhound
Greyhound
Greyhound
Greyhound
Greyhound

HT3

HT3

HT3

6MB L3 Cache

MPI task

Greyhound
Greyhound
Greyhound
Greyhound
Greyhound
Greyhound

6MB L3 Cache

MPI task

Greyhound
Greyhound
Greyhound
Greyhound
Greyhound
Greyhound

DDR3 Channel

DDR3 Channel

DDR3 Channel

DDR3 Channel

*To Interconnect*

HT1 / HT3

CRAY ✳

- 2 Multi-Chip Modules, 4 Opteron Dies

- 8 Channels of DDR3 Bandwidth to 8 DIMMs

- 24 (or 16) Computational Cores

  - 64 KB L1 and 512 KB L2 caches for each core

  - 6 MB of shared L3 cache on each die

- Dies are fully connected with HT3

- Snoop Filter Feature Allows 4 Die SMP to scale well

Use OpenMP across all 6 cores
in the Socket

# Proposed Programming Paradigm for Hybrid Multi-core

- MPI or PGAS between nodes and/or sockets
- OpenMP, Pthreads or some other shared memory parallelism across a portion of the cores on the node
- Vectorization to utilize the SSE# or SIMD units on the cores

- Node structure will be the same for all system sizes
  - Exascale
  - Petascale
  - Terascale

  ## Will all use the same node

CRAY
THE SUPERCOMPUTER COMPANY

# http://hybridmulticore.com/

# Cray XT6

| Characteristics | |
|---|---|
| Number of Cores | 16 or 24 (MC) |
| Peak Performance MC-8 (2.4) | 153 Gflops/sec |
| Peak Performance MC-12 (2.2) | 211 Gflops/sec |
| Memory Size | 32 or 64 GB per node |
| Memory Bandwidth | 83.5 GB/sec |

6.4 GB/sec direct connect HyperTransport

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

AMD Opteron 64

AMD Opteron 64

83.5 GB/sec direct connect memory

Cray SeaStar2+ Interconnect

# XT6 or XE6 Node Details: 24-core Magny Cours



- 2 Multi-Chip Modules, 4 Opteron Dies
- 8 Channels of DDR3 Bandwidth to 8 DIMMs
- 24 (or 16) Computational Cores, 24 MB of L3 cache
- Dies are fully connected with HT3
- Snoop Filter Feature Allows 4 Die SMP to scale well

# Cray SeaStar2+ Interconnect

**Now Scaled to 225,000 cores**

- **Cray XT6 systems ship with the SeaStar2+ interconnect**
- Custom ASIC
- **Integrated NIC / Router**
- MPI offload engine
- Connectionless Protocol
- Link Level Reliability
- Proven scalability to 225,000 cores

AMD Opteron 64

AMD Opteron 64

**6-Port Router**

DMA Engine

HyperTransport Interface

CRAY

Memory

Blade Control Processor Interface

PowerPC 440 Processor

# Why Custom Interconnects ?

- We have been keeping a close eye on progress with infiniband but we still believe custom interconnects are necessary for HPC



- What follows is a customer performance example showing the difference

- And a usage pattern we see on NSF systems…

# SeaStar and Infiniband Compared

- SeaStar-based XT5 systems are the first systems to scale application codes to the Petaflop level

- SeaStar provides scalability advantages at much smaller processor counts, *particularly on irregular communication patterns*

# Nearest Neighbor MPI Benchmark (Bucket Brigade, Large Messages)

*Both Infiniband and SeaStar scale well on a nearest neighbor test*

Legend:
- Cray XT4 (SeaStar)
- Tri Labs Opteron Cluster (Infiniband)

Y-axis: Mbytes/sec per MPI Rank
X-axis: Number of MPI Ranks (64, 256, 1024)

Data from *Red Storm / Cray XT4: A Superior Architecture for Scalability* by **Mahesh Rajan, Doug Doerfler, Courtenay Vaughan, Sandia National Laboratory -** **Presented at Cray User Group, May 4-9, 2009**

# Random Message MPI Benchmark Test (Short Messages)



This benchmark highlights the importance of a connectionless protocol in a scalable interconnect

Data from *Red Storm / Cray XT4: A Superior Architecture for Scalability* by **Mahesh Rajan, Doug Doerfler, Courtenay Vaughan, Sandia National Laboratory - Presented at Cray User Group, May 4-9, 2009**

# SIERRA/Presto – Weak Scaling

➤ Explicit Lagrangian mechanics with contact

➤ Model: Two sets of brick-walls colliding

➤ Weak scaling analysis with 80 bricks /PE, each discretized with 4x4x8 elements

➤ Contact algorithm communications dominates the run time

➤ *The rapid increase in run time after 64 processors on TLCC can be directly related to the poor performance on TLCC for random small-to-medium size messages*

➤ TLCC/Quad run time ratio at 1024 is 4X.

### PRESTO: Walls Collision Weak Scaling
### 10,240 Eelemnts/task; 596 Time Steps

# Usage Pattern – UT's Kraken Machine

| | |
|---|---|
| Award(U. Tennessee/ORNL) | Sep, 2007 |
| Cray XT3: 7K cores, 40 TF | Jun, 2008 |
| Cray XT4: 18K cores,166 TF | Aug 18, 2008 |
| Cray XT5: 65K cores, 600 TF | Feb 2, 2009 |
| Cray XT5+: ~100K cores, 1 PF | Oct, 2009 |

Kraken and Krakettes!

## XT5 CPU Usage by Core-Count

- 0 - 511 cores
- 512 - 1023 cores
- 1024 - 2047 cores
- 2048 - 4095 cores
- 4096 - 8191 cores
- 8192 - 16383 cores
- 16384 - 32767 cores
- 32768 - 66048 cores

5%
7%
14%
45%
10%
8%
9%
2%

Proprietary Cray Interconnect, SeaStar2, provides excellent scaling, with numerous tightly coupled applications running at 32K and 64K cores on the XT5.

NICS is specializing on true capability applications, plus high performance file and archival systems.

*Typical job size on IB cluster at TACC is ~300 cores*

# XT6 Processor Choices (vs. XT5)



| Processor | Frequency | Peak (Gflops) | Bandwidth (GB/sec) | Balance (bytes/flop) |
|---|---|---|---|---|
| Istanbul (XT5) | 2.6 | 62.4 | 12.8 | 0.21 |
| MC-8 | 2.0 | 64 | 42.6 | 0.67 |
| | 2.3 | 73.6 | 42.6 | 0.58 |
| | 2.4 | 76.8 | 42.6 | 0.55 |
| MC-12 | 1.9 | 91.2 | 42.6 | 0.47 |
| | 2.1 | 100.8 | 42.6 | 0.42 |
| | 2.2 | 105.6 | 42.6 | 0.40 |

# Cray XT6 Compute Blade

- New compute blade with 8 AMD Magny Cours processors
- Plug-compatible with XT5 cabinets and backplanes
- Initially will ship with SeaStar interconnect as the Cray XT6
- Upgradeable to Gemini Interconnect or Cray XE6
- Upgradeable to AMD's "Interlagos" series
- XT6 systems will continue to ship with the current SIO blade
- First customer ship, March 31st

# Cray XE6

# Cray XE6 Compute Blade

- **8 Magny Cours Sockets**
- **96 Compute Cores**
- **32 DDR3 Memory DIMMS**
- **32 DDR3 Memory channels**
- **2 Gemini ASICs**
- **L0 Blade management processor**

# Cray XE6 Compute Node

| Node Characteristics | |
|---|---|
| Number of Cores | 24 (Magny Cours) |
| Peak Performance MC-12 (2.2) | 211 Gflops/sec |
| Peak Performance MC-8 (2.4) | 153 Gflops/sec |
| Memory Size | 32 GB per node 64 GB per node |
| Memory Bandwidth (Peak) | 83.5 GB/sec |

# XE6 Node Details:
## 24-core Magny Cours



- 2 Multi-Chip Modules, 4 Opteron Dies
- 8 Channels of DDR3 Bandwidth to 8 DIMMs
- 24 (or 16) Computational Cores
  - 64 KB L1 and 512 KB L2 caches for each core
  - 6 MB of shared L3 cache on each die
- Dies are fully connected with HT3
- Snoop Filter Feature Allows 4 Die SMP to scale well

# Gemini Interconnect

# Cray Network Evolution

## SeaStar

- Built for scalability to 250K+ cores
- Very effective routing and low contention switch

## Gemini

- 100x improvement in message throughput
- 3x improvement in latency
- PGAS Support, Global Address Space
- Scalability to 1M+ cores

## Aries

- Ask me about it

# Cray Gemini

- 3D Torus network
- Supports 2 Nodes per ASIC
- 168 GB/sec routing capacity
- Scales to over 100,000 network endpoints
  - Link Level Reliability and Adaptive Routing
  - Advanced Resiliency Features
- Provides global address space
- Advanced NIC designed to efficiently support
  - MPI
    - Millions of messages/second
  - One-sided MPI
  - UPC, FORTRAN 2008 with coarrays, shmem
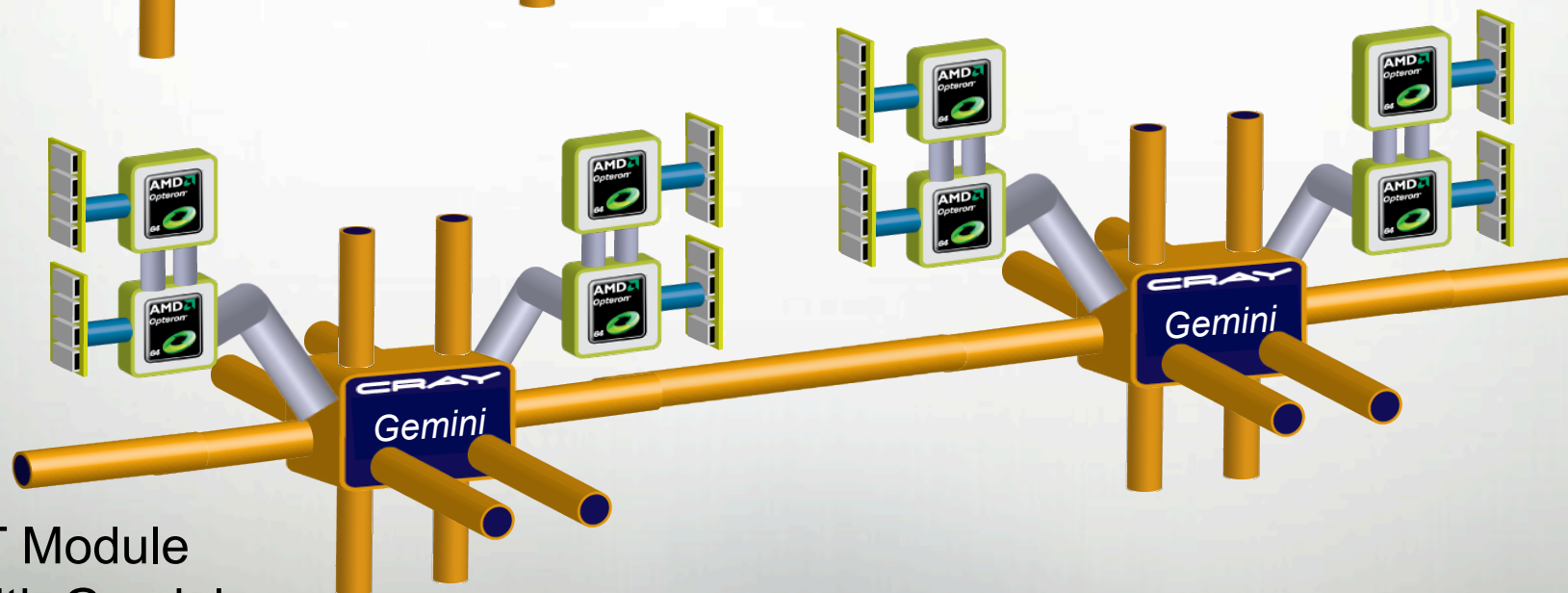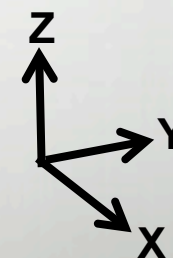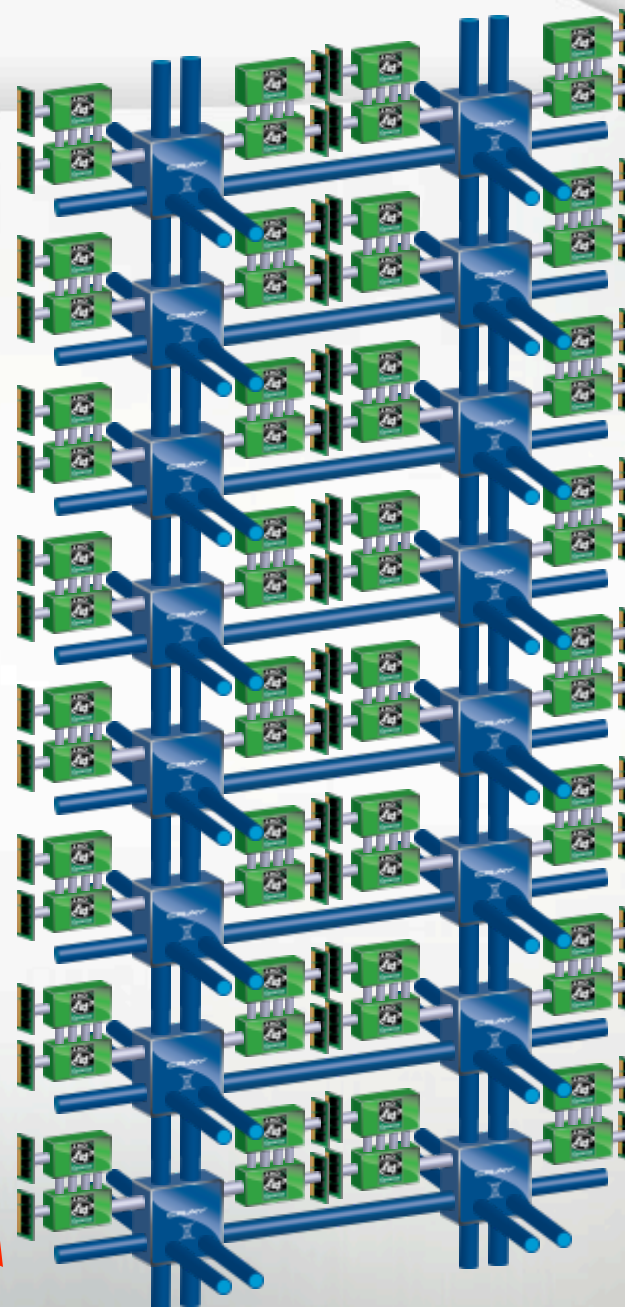  - Global Atomics



Hyper Transport 3

NIC 0

Hyper Transport 3

NIC 1

LO Processor

SB

Netlink

48-Port YARC Router

XT Module
with SeaStar

XT Module
with Gemini

Gemini

Gemini

Z

Y

X
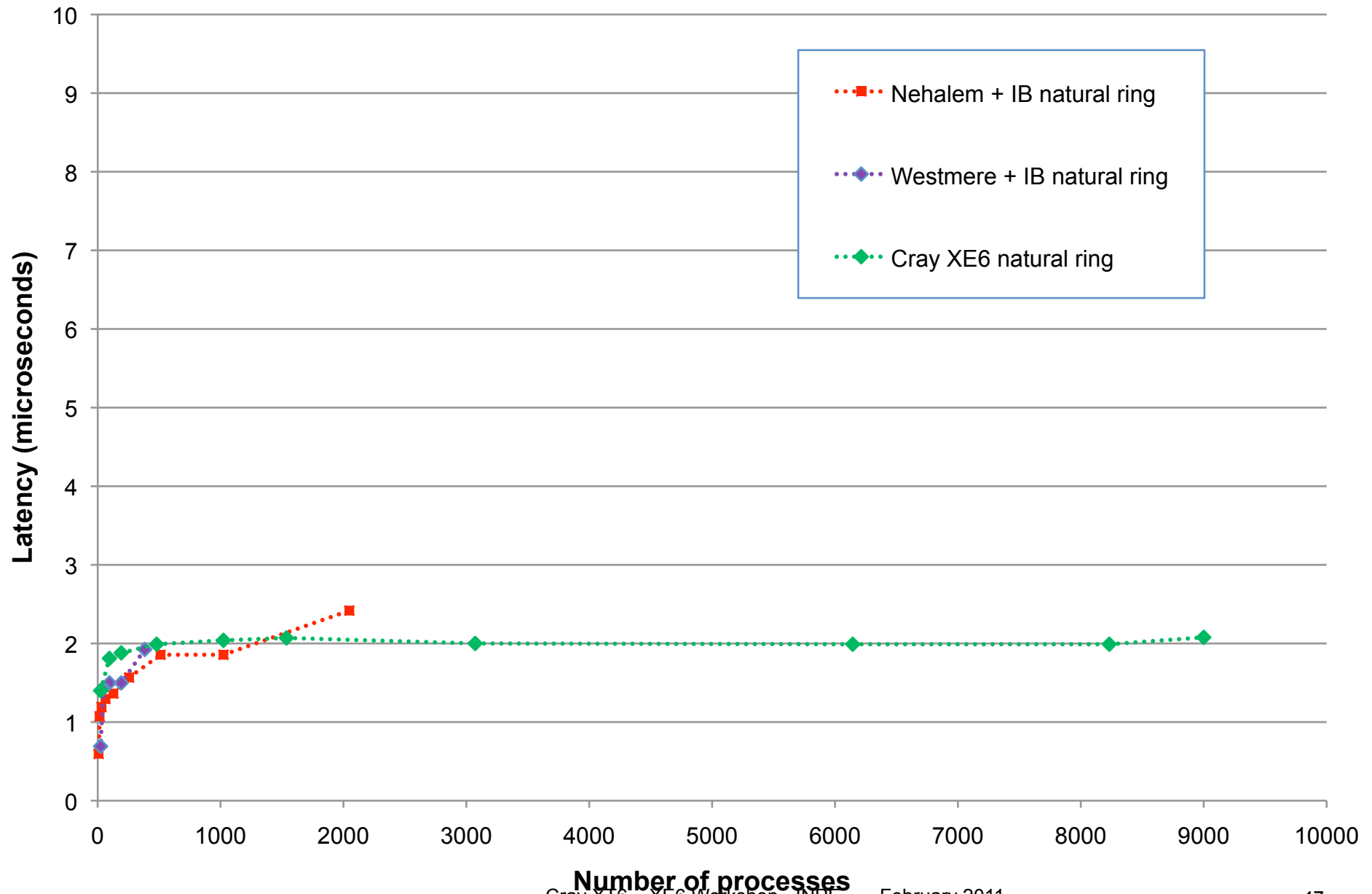
# Cray XE6 Chassis Topology

# Gemini Advanced Features

- Globally addressable memory provides efficient support for UPC, Co-array FORTRAN, Shmem and Global Arrays
  - Cray Programming Environment will target this capability directly

- Pipelined global loads and stores
  - Allows for fast irregular communication patterns

- Atomic memory operations
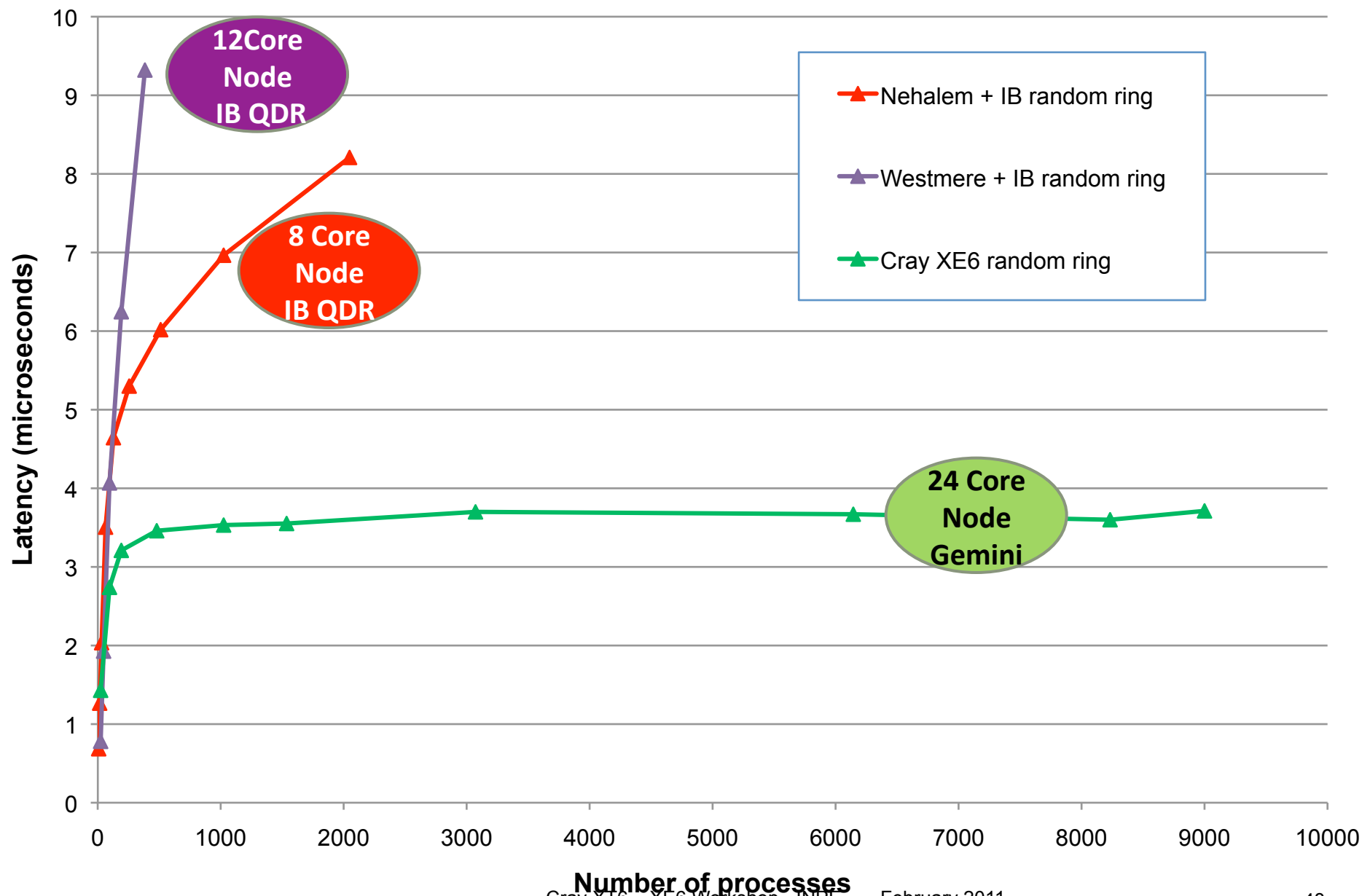  - Provides fast synchronization needed for one-sided communication models

# Gemini – QDR Comparison
## HPCC Natural Ring Latency Benchmark
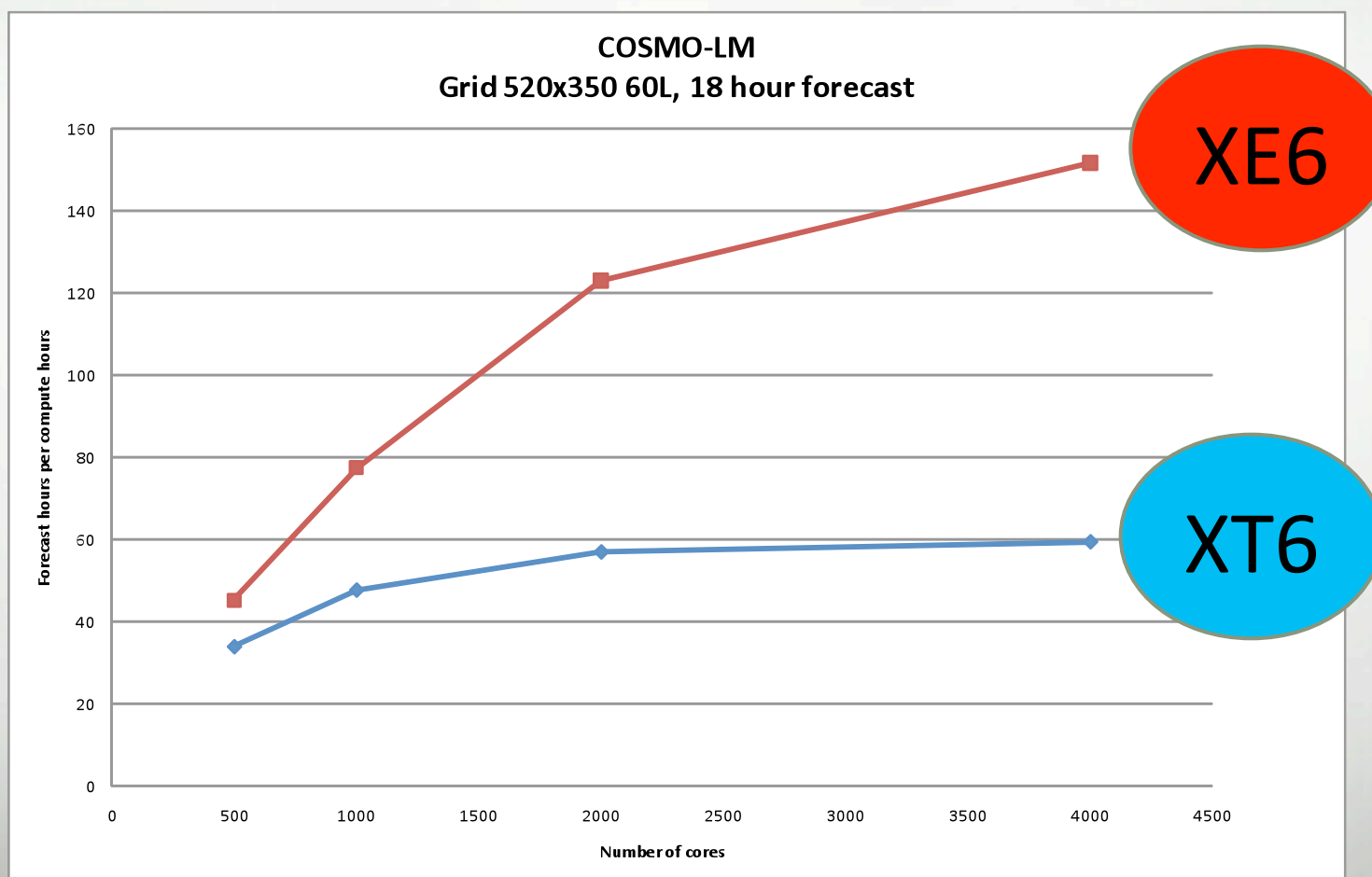
# Gemini – QDR Comparison
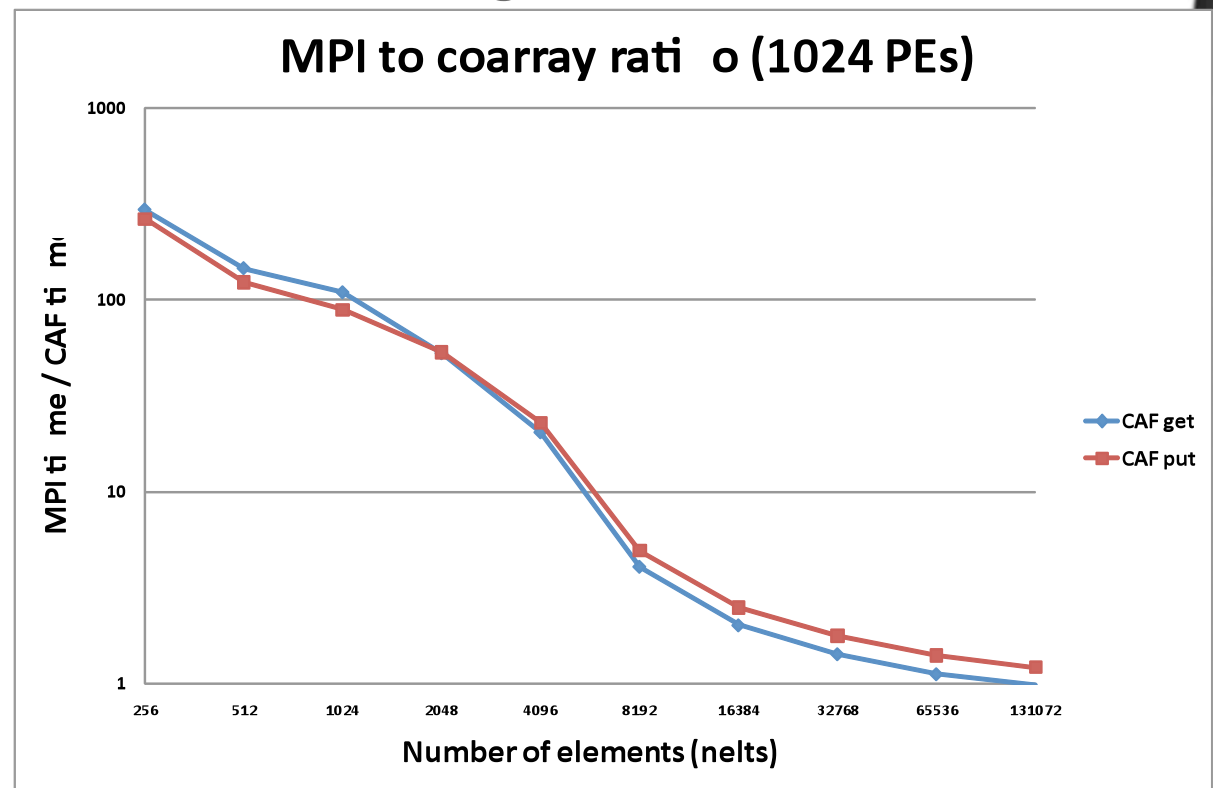## HPCC Random Ring Latency Benchmark

# Scalability and simulation rate

- Forecast Hours per compute Hours
- Typical performance improvement

**COSMO-LM**
**Grid 520x350 60L, 18 hour forecast**



XE6

XT6

# Remote gather: coarray vs MPI

- Coarray implementations are much simpler
- Coarray syntax allows the expression of remote data in a natural way – no need of complex protocols
- Coarray implementation is orders of magnitude faster for small numbers of indices

**MPI to coarray ratio (1024 PEs)**

*MPI time / CAF time* (y-axis)

Legend:
- CAF get
- CAF put

Number of elements (nelts): 256, 512, 1024, 2048, 4096, 8192, 16384, 32768, 65536, 131072

Thank You!